

CERTIFIED COPY OF  
PRIORITY DOCUMENT

日 本 国 特 許 庁  
JAPAN PATENT OFFICE

11036 U.S. PTO  
10/002762  
10/24/01

別紙添付の書類に記載されている事項は下記の出願書類に記載されて  
いる事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed  
with this Office

出 願 年 月 日

Date of Application:

2000年10月24日

出 願 番 号

Application Number:

特願2000-323823

出 願 人

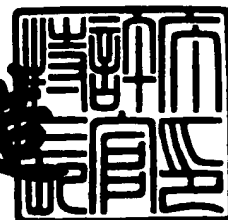
Applicant(s):

インターナショナル・ビジネス・マシーンズ・コーポレーシ  
ョン

2001年 6月19日

特 許 庁 長 官  
Commissioner,  
Japan Patent Office

及 川 耕 造



出証番号 出証特2001-3057749

【書類名】 特許願

【整理番号】 JP9000274

【提出日】 平成12年10月24日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 5/00

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 東京基礎研究所内

【氏名】 近藤 豪

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 東京基礎研究所内

【氏名】 広瀬 紳一

【特許出願人】

【識別番号】 390009531

【氏名又は名称】 インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】 100086243

【弁理士】

【氏名又は名称】 坂口 博

【代理人】

【識別番号】 100091568

【弁理士】

【氏名又は名称】 市位 嘉宏

【代理人】

【識別番号】 100106699

【弁理士】

【氏名又は名称】 渡部 弘道

【復代理人】

【識別番号】 100104880

【弁理士】

【氏名又は名称】 古部 次郎

【選任した復代理人】

【識別番号】 100100077

【弁理士】

【氏名又は名称】 大場 充

【手数料の表示】

【予納台帳番号】 081504

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9706050

【包括委任状番号】 9704733

【包括委任状番号】 0004480

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 構造回復システム、構文解析システム、変換システム、コンピュータ装置、構文解析方法、記憶媒体及びプログラム伝送装置

【特許請求の範囲】

【請求項 1】 所定の規則に基づいて記述されたデータ列の構造を解析し、当該規則上の誤りを検出する解析手段と、

前記解析手段からの依頼に応じて、前記解析手段により検出された前記データ列における前記規則上の誤りを修正する回復手段とを備え、

前記回復手段は、

特定の種類の誤りを修正する単純な機能を有する修正手段の集合を含んで構成され、

前記データ列における前記規則上の誤りの種類に応じて前記修正手段を切り換えて使用することにより、前記データ列における種々の誤りを修正することを特徴とする構造回復システム。

【請求項 2】 所定の規則に基づいて記述されたデータ列に対する構文解析を行う構文解析システムにおいて、

構文解析処理を実行する構文解析部と、

前記構文解析部からの依頼に応じて、前記構文解析部の構文解析処理により検出された前記データ列の誤りを修正する構文回復部とを備え、

前記構文回復部は、前記修正の内容を変更可能であることを特徴とする構文解析システム。

【請求項 3】 前記構文回復部は、前記構文解析部により検出される前記データ列の誤りの種類に応じて複数用意され、

個々の前記構文回復部は、特定の 1 種類の誤りを修正する一つの機能を有すること

を特徴とする請求項 2 に記載の構文解析システム。

【請求項 4】 データ列の種類と、当該データ列の誤りを回復するために用いる前記構文回復部の組み合わせとを対応付けた対応情報を保存する対応情報保存手段をさらに備え、

前記構文解析部は、処理対象である前記データ列の種類に応じて、前記対応情報保存手段に保存された前記対応情報に基づいて、誤りの修正を依頼する前記構文回復部の組み合わせを設定すること

を特徴とする請求項 3 に記載の構文解析システム。

【請求項 5】 前記構文回復部の少なくとも一つは、

前記データ列が前記構文解析部における構文解析処理に用いる規則において定義されていない要素を含んでいる場合に起動し、

前記構文解析部にて用いられる規則を前記要素を定義している規則に置き換え、前記データ列を前記構文解析部に戻す処理を行うこと

を特徴とする請求項 3 に記載の構文解析システム。

【請求項 6】 処理対象である前記データ列に対する字句の解析処理を実行する字句解析部と、

前記字句解析部からの依頼に応じて、前記字句解析部の解析処理により検出された前記データ列における字句の誤りを修正する字句回復部とをさらに備え、

前記字句回復部は、前記修正の内容を変更可能であること  
を特徴とする請求項 2 に記載の構文解析システム。

【請求項 7】 前記字句回復部は、前記字句解析部により検出される前記データ列における字句の誤りの種類に応じて複数用意され、

個々の前記字句回復部は、特定の 1 種類の誤りを修正する一つの機能を有すること

を特徴とする請求項 6 に記載の構文解析システム。

【請求項 8】 所定の形式で記述されたデータ列を他の形式に変換する変換システムにおいて、

前記データ列を解析する解析部と、

前記解析部からの依頼に応じて、前記解析部の解析処理により検出された前記データ列の誤りを修正する回復部と、

前記解析部の解析結果に基づいてデータ形式の変換を行う変換部とを備え、

前記回復部は、前記解析部により検出される前記データ列の誤りの種類に応じて複数用意され、

個々の前記回復部は、特定の 1 種類の誤りを修正する一つの機能を有することを特徴とする変換システム。

【請求項 9】 前記解析部は、前記データ列に対する構文解析を行う構文解析手段であり、

前記回復部は、前記データ列における構文規則上の誤りを修正する構文回復手段であること

を特徴とする請求項 8 に記載の変換システム。

【請求項 10】 所定の規則に基づいて記述されたデータ列を入力する入力部と、プログラム制御により実現される機能によって当該データ列を処理する演算処理部と、当該演算処理部により処理された前記データ列を出力する出力部とを備えたコンピュータ装置において、

前記演算処理部は、

前記データ列を解析する解析部と、

前記解析部からの依頼に応じて、前記解析部の解析処理により検出された前記データ列の誤りを修正する回復部とを備え、

前記回復部は、前記解析部により検出される前記データ列の誤りの種類に応じて複数用意され、

個々の前記回復部は、特定の 1 種類の誤りを修正する一つの機能を有することを特徴とするコンピュータ装置。

【請求項 11】 前記解析部は、前記データ列に対する構文解析を行う構文解析手段であり、

前記回復部は、前記データ列における構文規則上の誤りを修正する構文回復手段であること

を特徴とする請求項 10 に記載のコンピュータ装置。

【請求項 12】 所定の規則に基づいて記述されたデータ列に対する構文解析を行う構文解析方法において、

処理対象である前記データ列に構文規則上の誤りがあった場合に、当該誤りを修正するために使用するプログラムモジュールを選択するステップと、

前記データ列に対する構文解析を行うステップと、

構文解析において前記データ列から構文規則上の誤りが検出された場合に、前記プログラムモジュールに修正を依頼するステップと、

前記プログラムモジュールによる修正を行った後、修正された前記データ列に対する構文解析を行うステップとを含むことを特徴とする構文解析方法。

【請求項 1 3】 前記使用するプログラムモジュールを選択するステップは、

処理対象である前記データ列の種類を調べるステップと、

得られた前記データ列の種類に応じて、予め定められた対応関係に基づいて使用する前記プログラムモジュールを選択するステップと

を含むことを特徴とする請求項 1 2 に記載の構文解析方法。

【請求項 1 4】 前記プログラムモジュールに修正を依頼した後、当該プログラムモジュールによる指示に従って、前記データ列に対する構文解析に用いる規則を他の規則に置き換えるステップをさらに備え、

修正後のデータ列に対する構文解析を行うステップにおいて、置き換えられた前記他の規則に基づいて前記データ列に対する構文解析を行うことを特徴とする請求項 1 2 に記載の構文解析方法。

【請求項 1 5】 コンピュータに実行させるプログラムを当該コンピュータの入力手段が読取可能に記憶した記憶媒体において、

前記コンピュータに実現させる機能として

データ列を解析する解析手段と、

前記解析手段からの依頼に応じて、前記解析手段の解析処理により検出された前記データ列の誤りを修正する回復手段とを備え、

前記回復手段は、前記解析処理により検出される前記データ列の誤りの種類に応じて複数用意され、

個々の前記回復手段は、特定の 1 種類の誤りを修正する一つの機能を有するプログラム・プロダクトを格納したことを特徴とする記憶媒体。

【請求項 1 6】 コンピュータに実現させる機能として、データ列を解析する解析手段と、前記解析手段からの依頼に応じて、前記解析手段の解析処理によ

り検出された前記データ列の誤りを修正する回復手段とを備え、前記回復手段は、前記解析処理により検出される前記データ列の誤りの種類に応じて複数用意され、個々の前記回復手段は、特定の 1 種類の誤りを修正する一つの機能を有するプログラム・プロダクトを記憶する記憶手段と、

前記記憶手段から前記プログラム・プロダクトを読み出して当該プログラムを送信する送信手段とを備えたことを特徴とするプログラム伝送装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、所定の規則に基づいて記述されたプログラムや文書などの字句解析及び構文解析を行い、誤りを修正する技術に関する。

【0002】

【従来の技術】

所定の規則に基づいて記述されたプログラムや文書の誤りを修正するための技術として、当該プログラムや文書などのデータ列に対して字句解析や構文解析を行い、誤りを検出する技術がある。従来のこの種の技術では、主に、次の 2 種類の方法で誤りを扱っていた。

【0003】

第 1 の方法は、誤りを検出した場合に警告を発し、解析を中断するか、またはその誤りの箇所よりも後方にある同期ポイントから解析を再開する方法である。すなわち、この方法では、積極的に誤りを正常な状態に回復することはしない。これは、プログラミング言語の処理系などのように解析対象の文書に誤りを許さない場合や、アプリケーションを特定しない汎用的な部品として開発された解析システムに多く採用されている手法である。

【0004】

例として、XML (eXtensible Markup Language) の処理系について説明する。XML では、アプリケーションに依存しない汎用的な構文解析システムが、数々のベンダーから提供されている（例えば、米国オープン XML 社の OpenXML や、米国 IBM 社の XML4J など）。これらの処理系では、文書中に誤りがあった場



合、解析自体を諦めてしまうか、誤りを起こした字句を無視して後方の同期ポイントから再開する処理を行っている。かかる処理を行う際に、構文解析システムを操作するインターフェース（SAX: Simple API for XML）に基づいて、ErrorHandlerという外部モジュールが、誤りの起こったポイントと誤りの説明的なメッセージを受け取ることができる。しかし、このモジュールは、単に警告を受け取っているのみであり、構文解析システムの状態や出力結果を変化させる機能は提供されていない。

## 【0005】

第2の方法は、解析システムを用いているアプリケーションに対応した固有の回復方法に基づいて解析結果を出力し、解析を続行する方法である。すなわち、この方法では、単に誤りを検出するのみでなく、積極的に誤りを修正し、構文規則などの誤りのないプログラムや文書を回復する。この方法は、ソース文書の作成者と閲覧者が全く異なっており、文書中に誤りがあっても閲覧者には何らかの出力結果を生成しなければならない場合に用いられる。

## 【0006】

例として、HTML (HyperText Markup Language) の処理系について説明する。インターネットなどで用いられるウェブページを記述するHTMLは、通常、HTML文書の作成者と閲覧者が異なる。そのため、HTML文書に構文規則などの誤りがあったとしても、閲覧者に対してその誤りの存在を提示するだけでは足りず、HTML文書としてのつじつまを合わせて閲覧可能な状態にする必要がある。したがって、ウェブブラウザ（ウェブページの閲覧用ソフトウェア）には、上記のような字句解析や構文解析を行う機能と、規則上の誤りが検出されたHTML文書を誤りのない状態に回復する機能とを有しているものがある。

## 【0007】

所定のHTML文書中に

```
<P>str0 <B>str1 <I>str2</B> str3</I> str4</P>
```

という断片があったとする。この断片は、タグ<B> </B>と<I> </I>とが入れ子構

造になっていないため、構文規則において誤りである。しかし、ウェブブラウザに表示を行うためには、この誤りをそのまま放置することはできず、ウェブブラウザの構文解析手段は何らかの出力を生成しなければならない。

#### 【 0 0 0 8 】

代表的なウェブブラウザである米国ネットスケープ・コミュニケーションズ社の Netscape Navigator は、上記のような断片に対して

```
<P>str0 <B>str1 <I>str2</I> str3</B> str4</P>
```

のように修正する。すなわち、</B>と</I>とを入れ替えることにより入れ子構造を回復する。これにより、図 1 4 に示すような出力結果（str1とstr3は太字、str2は太字で斜体）が得られる。

#### 【 0 0 0 9 】

これに対し、他の代表的なウェブブラウザである米国マイクロソフト社の Internet Explorer は、上記のような断片に対して

```
<P>str0 <B>str1 <I>str2</I></B><I> str3</I> str4</P>
```

のように修正する。すなわち、</B>の前に</I>を補い、</B>の後に<I>を補うことにより入れ子構造を回復する。これにより、図 1 5 に示すような出力結果（str1は太字、str2は太字で斜体、str3は斜体）が得られる。

#### 【 0 0 1 0 】

図 1 2 は、従来の構文解析システムの構成例を示す図、図 1 3 は、図 1 2 に示す構文解析システムにおける構文解析処理を説明するフローチャートである。

図 1 2 を参照すると、構文解析システム 1 2 0 は、入力文書中の一定の入力ストリームを入力し字句解析を行う字句解析部 1 2 1 と、字句解析部 1 2 1 から字句を取得して構文解析を行い入力文書の構造を示す抽象構文木（AST）を生成して出力する構文解析部 1 2 2 と、抽象構文木の生成に用いられるノード作成部 1 2 3 とを備える。字句解析部 1 2 1 は、字句解析に用いるバッファ 1 2 1 a と

、字句の誤りを修正する字句回復部 1 2 1 b とを備える。構文解析部 1 2 2 は、構文解析に用いるバッファ 1 2 2 a 及び文脈ポインタと、構文規則の誤りを修正する構文回復部 1 2 2 b とを備える。また、構文解析部 1 2 2 による処理が行われる際に、構文解析に用いられる文法情報オブジェクト 1 2 4 が生成される。

#### 【 0 0 1 1 】

図 1 3 に示すように、構文解析処理が開始されると、まず、構文解析部 1 2 2 が初期化される（ステップ 1 3 0 1）。構文解析部 1 2 2 の初期化は、次の三つの作業により行われる。すなわち、①入力文書に対応する文書の種類を解析し、文法情報オブジェクト 1 2 4 を生成する。②バッファ 1 2 2 a を空にする。③文脈ポインタが抽象構文木のルートノードを指すようにする。なお、構文解析部 1 2 2 の初期化に先立って、初期的に入カストリームの入力及び字句解析は済んでいるものとする。

#### 【 0 0 1 2 】

次に、構文解析部 1 2 2 が、バッファ 1 2 2 a から字句を取り出し、処理対象である字句 t とする（ステップ 1 3 0 2）。バッファ 1 2 2 a が空の場合は（ステップ 1 3 0 1 で初期化した直後は常に空の状態となっている）、字句解析部 1 2 1 に字句を要求し、取得した字句を字句 t とする。字句 t が入力文書の終端であった場合は、生成された抽象構文木を出力して処理を終了する（ステップ 1 3 0 3）。

#### 【 0 0 1 3 】

字句 t が入力文書の終端でない場合、構文解析部 1 2 2 は、字句 t が文脈ポインタに対して文法的に適合するかどうかを文法情報オブジェクト 1 2 4 に問い合わせる。そして、適合するならば文脈ポインタに字句 t を追加する（ステップ 1 3 0 4 で Y e s）。この追加は、次の処理によって行う。まず、非終端記号のノード n をノード作成部 1 2 3 により作成し、文脈ポインタに追加する（ステップ 1 3 0 5）。次に、文脈ポインタの指示先を新たに追加した非終端記号のノード n にずらす（ステップ 1 3 0 6）。そして、文脈ポインタが指し示す非終端記号のノード n が全ての子供のノードを持った場合、文脈ポインタを親ノードへずらす（ステップ 1 3 0 7、1 3 0 8）。文脈ポインタが指し示す非終端記号のノー

ド n が全ての子供のノードを持っていない場合、またはステップ 1 3 0 8 で文脈ポインタを親ノードへずらした後、ステップ 1 3 0 2 へ戻り、次の字句を取得して同様の処理を行う。

#### 【0 0 1 4】

ステップ 1 3 0 4 において、字句 t が文脈ポインタに対して文法的に適合しない場合、構文解析部 1 2 2 はエラーの内容を出力する（ステップ 1 3 0 9）。そして、予め定められたエラー処理を実行した後（ステップ 1 3 1 0）、ステップ 1 3 0 2 へ戻り、次の字句を取得して処理を行う。ここで、エラー処理とは、当該字句 t を読み飛ばして次の字句に対する処理に移行することや、固定的な手法による回復処理を含む。回復処理を行う場合、構文解析部 1 2 2 は、構文回復部 1 2 2 b を呼び出し、字句 t が文脈ポインタに対して文法的に適合するように修正した後にステップ 1 3 0 2 へ戻ることとなる。

#### 【0 0 1 5】

このような回復処理は、構文解析システム 1 2 0（ウェブブラウザの構文解析手段）が HTML 専用であり、かつアプリケーション専用の独自に開発されたシステムであるために行うことができた。かかる手法を用いた構文解析システムとしては、例えば、株式会社ジャストシステムの Ark や、W 3 C（World Wide Web Consortium）による W3C Tidy などがある。

なお、上記の動作例で、字句解析部 1 2 1 による字句解析の際に、必要に応じて字句回復部 1 2 1 b による字句の誤りの修正が行われるが、入力ストリームの字句を所定の規則に照らして適切な字句に置き換える単純な処理であるため、説明を省略した。

#### 【0 0 1 6】

##### 【発明が解決しようとする課題】

上述したように、構文解析システムが特定のプログラミング言語専用であり、所定のアプリケーションに固有のシステムであれば、構文規則などに誤りのあるプログラムや文書を修正し、誤りのない状態に回復することが可能である。

#### 【0 0 1 7】

しかし、上記の HTML 文書の場合のように、誤りのある文書を正常な状態に

回復させる手法が複数存在する場合があるのに対し、従来の構文解析システムは、特定の誤りに対して固定的な手法で修正を行っていた。そのため、必ずしも HTML 文書の作者の意図に応じた修正が行われるとは限らなかった。

## 【 0 0 1 8 】

さらに、文書の誤りを修正する手法は、解析システムの作成者が定義したものだけではなく、利用者自身が定義することも含め、様々な条件によって選択可能であることが望ましい。

## 【 0 0 1 9 】

そこで本発明は、プログラムや文書などのデータ列における字句や構文規則の誤りの種類に応じて複数の回復手段を用意し、選択的に適用して誤りを修正できるシステムを提供することを目的とする。

## 【 0 0 2 0 】

## 【課題を解決するための手段】

上記の目的を達成するため、本発明は、所定の規則に基づいて記述されたデータ列の構造を解析し、規則上の誤りを検出する解析手段と、この解析手段とは独立に設けられ、この解析手段からの依頼に応じて、この解析手段により検出されたこのデータ列における規則上の誤りを修正する回復手段とを備え、この回復手段は、特定の種類の誤りを修正する単純な機能を有する修正手段の集合を含んで構成され、このデータ列における規則上の誤りの種類に応じてこの修正手段を切り換えて使用することにより、このデータ列における種々の誤りを修正することを特徴とする。

これにより、処理対象であるデータ列の種類や誤りの種類に応じて修正手段を作成し、追加、変更または不要な修正手段を削除することにより、種々のデータ列や種々の誤りに対して柔軟に対応し、適切な修正を行うことが可能となる。

## 【 0 0 2 1 】

特に、本発明は、プログラムや文書における構文解析に用いることができる。この場合、本発明は、所定の規則に基づいて記述されたデータ列に対する構文解析を行う構文解析システムにおいて、構文解析処理を実行する構文解析部と、この構文解析部からの依頼に応じて、この構文解析部の構文解析処理により検出さ

れたこのデータ列の誤りを修正する構文回復部とを備え、この構文回復部は、構文解析部とは独立に設けられ、かつ修正の内容を変更可能であることを特徴とする。

#### 【 0 0 2 2 】

ここで、この構文回復部は、構文解析部により検出されるデータ列の誤りの種類に応じて複数用意される。そして、個々の構文回復部は、特定の 1 種類の誤りを修正する一つの機能を有する。

したがって、この構文回復部を追加、変更、削除することにより、柔軟かつ適切な誤りの修正を実現することができる。

#### 【 0 0 2 3 】

ここで、適切な構文回復部を用いるための手段として、データ列の種類と、このデータ列の誤りを回復するために用いる構文回復部の組み合わせとを対応付けた対応情報を保存する対応情報保存手段をさらに備える構成とし、構文解析部は、処理対象であるデータ列の種類に応じて、この対応情報保存手段に保存された対応情報に基づいて、誤りの修正を依頼する構文回復部の組み合わせを設定することができる。

この対応情報保存手段としては、データ列の種類を識別するための識別子（ID）と、対応する構文回復部の組を登録した対応表を用いることができる。

#### 【 0 0 2 4 】

さらに、この構文回復部の少なくとも一つは、処理対象のデータ列が構文解析部における構文解析処理に用いる規則において定義されていない要素を含んでいる場合に起動し、構文解析部にて用いられる規則を、このデータ列に含まれている当該要素を定義している規則に置き換え、このデータ列を構文解析部に戻す処理を行うことを特徴とする。

#### 【 0 0 2 5 】

さらにまた、この構文解析システムは、字句解析においても同様の構成を取ることができる。すなわち、処理対象である前記データ列に対する字句の解析処理を実行する字句解析部と、この字句解析部からの依頼に応じて、この字句解析部の解析処理により検出されたこのデータ列における字句の誤りを修正する字句回

復部とをさらに備え、この字句回復部は、字句解析部とは独立に設けられ、かつ修正の内容を変更可能であることを特徴とする。

そして、この字句回復部は、字句解析部により検出されるデータ列における字句の誤りの種類に応じて複数用意される。また、個々の字句回復部は、特定の 1 種類の誤りを修正する一つの機能を有する。

これにより、字句の誤りの修正においても柔軟かつ適切な処理を行うことができる。

#### 【 0 0 2 6 】

また、本発明は、上記のような構文解析システムを含むデータ変換システムとして提供することもできる。

さらに、上記のような構文解析システムやデータ変換システムを実現するコンピュータ装置として提供することもできる。

さらにまた、本発明は、上記のような構文解析システムをコンピュータにより実現させるプログラム・プロダクトとして作成し、かかるプログラム・プロダクトを記憶した記憶媒体や、ネットワークを介してこのプログラム・プロダクトを伝送する伝送装置として提供することができる。

#### 【 0 0 2 7 】

また、本発明は、所定の規則に基づいて記述されたデータ列に対する構文解析を行う構文解析方法において、処理対象であるデータ列に構文規則上の誤りがあった場合に、この誤りを修正するために使用するプログラムモジュールを選択するステップと、このデータ列に対する構文解析を行うステップと、構文解析においてこのデータ列から構文規則上の誤りが検出された場合に、プログラムモジュールに修正を依頼するステップと、このプログラムモジュールによる修正を行った後、修正されたデータ列に対する構文解析を行うステップとを含むことを特徴とする。

#### 【 0 0 2 8 】

ここで、この使用するプログラムモジュールを選択するステップは、処理対象であるデータ列の種類を調べるステップと、得られたデータ列の種類に応じて、予め定められた対応関係に基づいて使用するプログラムモジュールを選択するス

テップとを含む構成とすることができる。

【0029】

さらに、この構文解析方法は、プログラムモジュールに修正を依頼した後、このプログラムモジュールによる指示に従って、構文解析に用いる規則を他の規則に置き換えるステップをさらに備え、修正後のデータ列に対する構文解析を行うステップにおいて、置き換えられた他の規則に基づいて処理対象のデータ列に対する構文解析を行う構成とすることができる。

【0030】

【発明の実施の形態】

以下、添付図面に示す実施の形態に基づいて、この発明を詳細に説明する。

まず、本発明の概要について説明する。本発明は、プログラムや文書などの所定の規則に基づいて記述されたデータ列を他の形式に変換する場合などに用いられる構文解析システムにおいて、字句や構文規則の誤りを修正する手段を字句解析手段及び構文解析手段（以下、これらを特に区別しない場合は単に解析手段と称す）から独立させ、誤りの種類ごとの単純な修正機能を持つ回復手段群を用意する。そして、字句解析手段及び構文解析手段が誤りを発見した場合は、これら回復手段群に回復処理を依頼する。

【0031】

処理を依頼された回復手段は、誤りの種類及び誤りを起こした字句または文字列に基づいて解析手段の出力結果及びバッファを操作する。これらの回復手段群は、解析手段に対して独立しているので、任意に追加、削除が可能である。

【0032】

また、これらの回復手段群が行う回復処理は、お互いに独立でない場合が考えられる。すなわち、所定の文書に対して所定の回復手段の列（いくつかの回復手段の組み合わせ）により回復可能であったが、別の文書に対しては当該回復手段の列が有効ではなく、別の列を用いる必要があるといった場合である。このような場合に対処するために、文書の識別子と、それを解析するために必要な誤り回復手段の列の対応関係を保存する。そして、所定の文書に対する解析処理の実行時に、保存されている対応関係に基づいて、処理対象である当該文書に対応する



回復手段の列をシステムに設定する。

### 【 0 0 3 3 】

図 1 は、本発明の実施の形態における構文解析システムの全体構成を示す図である。

図 1 を参照すると、本実施の形態の構文解析システム 1 0 0 は、入力文書中の一定の入力ストリームを入力し字句解析を行う字句解析部 1 0 と、字句の誤りを修正する字句回復部 1 2 と、字句解析部 1 0 から字句を取得して構文解析を行い入力文書の構造を示す抽象構文木 ( A S T ) を生成して出力する構文解析部 2 0 と、構文規則の誤りを修正する構文回復部 2 2 と、使用する構文回復部 2 2 を指定するために用いる構文回復部 2 2 の対応表 3 0 と、抽象構文木の生成に用いられるノード作成部 2 3 とを備える。また、構文解析部 2 0 による処理が行われる際に、構文解析に用いられる文法情報オブジェクト 2 4 が生成される。さらに、字句解析部 1 0 は、字句解析処理に用いるバッファ 1 1 を備え、構文解析部 2 0 は、構文解析処理に用いるバッファ 2 1 及び文脈ポインタを備える。

### 【 0 0 3 4 】

図 1 において、字句解析部 1 0 、字句回復部 1 2 、構文解析部 2 0 、構文回復部 2 2 、ノード作成部 2 3 及び文法情報オブジェクト 2 4 は、コンピュータプログラムにより制御された C P U にて実現される仮想的なソフトウェアブロックである。特に、字句回復部 1 2 及び構文回復部 2 2 は、プログラムモジュールとして構文解析システム 1 0 0 の使用者が任意に作成し、追加、変更、削除することができる。C P U を制御する当該コンピュータプログラムは C D - R O M やフロッピーディスクなどの記憶媒体に格納したり、ネットワークを介して伝送したりすることにより提供される。

### 【 0 0 3 5 】

図 1 1 は、本実施の形態の構文解析システム 1 0 0 を搭載したコンピュータ装置の構成例を説明する図である。

図 1 1 を参照すると、コンピュータ装置 2 0 0 は、C P U 2 1 0 と、C P U 2 1 0 を制御して図 1 に示した構文解析システム 1 0 0 の各構成要素を実現するコンピュータプログラムを格納したメモリ 2 2 0 と、処理対象である入力文書を入

力する入力部 2 3 0 と、入力文書に対する構文解析の結果を出力する出力部 2 4 0 とを備える。また、コンピュータ装置 2 0 0 は、構文解析システム 1 0 0 を実現するコンピュータプログラムを、上述した C D - R O M やフロッピーディスクから読み出すディスクドライブ 2 5 0 や、ネットワークを介して受信する受信部 2 6 0 を備える。

このコンピュータ装置 2 0 0 は、例えば、H T M L 文書を W M L (Wireless Markup Language) 文書に変換するというような、文書の形式を変換する変換システム (コンバータ) として使用され、文書形式の変換に先立って元の文書の構文解析を行う際に、構文解析システム 1 0 0 を動作させる。また、コンピュータ装置 2 0 0 をこのような変換システムとして用いる場合は、上述した出力部 2 4 0 の出力結果 (構文解析システム 1 0 0 は、本実施の形態で説明するように、入力文書の字句及び構文規則の誤りを修正するので、この出力結果は、これらの誤りが修正されたものである) に基づいて、当該入力文書の形式を変換する変換部 (図示せず) をさらに備えることとなる。

#### 【 0 0 3 6 】

本実施の形態の構文解析システム 1 0 0 は、入力文書における構文規則の誤りを修正する手段として、構文回復部 2 2 を、構文解析部 2 0 から独立して設けている。また、構文回復部 2 2 は、特定の構文規則の誤りを修正する単純な機能を有しており、構文解析処理の際には、複数の構文回復部 2 2 を組み合わせて用いることにより、種々の構文規則の誤りに対応する。

#### 【 0 0 3 7 】

構文回復部 2 2 の組み合わせは、処理対象である入力文書の種類に応じて決定される。入力文書の種類に応じた構文回復部 2 2 の組み合わせの設定は、予め適当な基準を設けて行っても良いし、実際に構文解析を行った経験に基づいて行っても良い。また、入力文書の種類の分類基準は、構文解析システム 1 0 0 の使用者が任意に決定することができる。

例えば、S G M L などのマークアップ言語を入力文書とする場合、S G M L において文書の形式を指定する D T D (Document Type Definition: 文書型定義) に基づいて分類することもできる。

さらに詳細な分類を設定することも可能である。例として、HTML文書（ウェブページ）を入力文書とする場合であって、特定のサイトにおいて、サイト内のHTML文書に対して特定の処理を行っている場合を考える。そして、HTML文書に対して処理を行うプログラムの性質上、HTML文書に特定の構文規則の誤りを発生させやすいものとする。この場合、当該サイト内に存在するHTML文書は、同様の構文規則の誤りが多いというような特定の傾向がある場合が考えられる。そこで、このサイト内のHTML文書を上述した入力文書の種類の分類項目としておき、これに対応する構文回復部22の組み合わせを設定することが可能である。

設定された構文回復部22の組み合わせは、入力文書の種類を示す文書IDと対応付けて、対応表30に登録される。そして、構文解析処理を行う際に、入力文書の種類に応じた文書IDをキーとして検索することにより、適切な構文回復部22の組み合わせを特定することができる。

#### 【0038】

本実施の形態では、字句回復部12と字句解析部10との関係も、構文回復部22と構文解析部20との関係と同様である。すなわち、字句回復部12を、字句解析部10から独立して設けてある。また、字句回復部12は、特定の字句の誤りを修正する単純な機能を有し、字句解析処理の際には、複数の字句回復部12を組み合わせる用いることにより、種々の字句の誤りに対応する。

ただし、字句の誤りは、ある程度、種類が限定されるため、構文回復部22の場合のように入力文書の種類に応じた組み合わせを用意せず、入力ストリームに対して常に全ての字句回復部12を適用するようにしても現実的な処理が可能である。したがって、本実施の形態では、対応表30のような入力文書の種類に対する対応情報を保存する手段は設けていない。しかしながら、字句回復部12に関しても、構文回復部22と同様に、入力文書の種類に応じた対応情報を保存する手段を設け、入力文書の種類ごとに適切な字句回復部12の組み合わせを特定して処理を行うこともできるのは言うまでもない。

#### 【0039】

次に、本実施の形態の動作について説明する。

図 2 は、本実施の形態の構文解析部 2 0 による構文解析処理の流れを説明するフローチャートである。

図 2 に示すように、構文解析処理が開始されると、まず、構文解析部 2 0 が初期化される（ステップ 2 0 1）。構文解析部 2 0 の初期化は、従来と同様に、①文法情報オブジェクト 2 4 の生成、②バッファ 2 1 を空にすること、③文脈ポインタが抽象構文木のルートノードを指すようにすること、の三つの作業を行い、さらに、対応表 3 0 を参照して、入力文書の種類に基づき、対応する構文回復部 2 2 の組を獲得し、該当する構文回復部 2 2 を構文解析処理において使用するよう設定する。

なお、構文解析部 2 0 の初期化に先立って、初期的に入力ストリームの入力及び字句解析は済んでいるものとする。また、字句解析部 1 0 の動作については後述する。

#### 【 0 0 4 0 】

次に、構文解析部 2 0 が、バッファ 2 1 から字句を取り出し、処理対象である字句  $t$  とする（ステップ 2 0 2）。バッファ 2 1 が空の場合は（ステップ 2 0 1 で初期化した直後は常に空の状態となっている）、字句解析部 1 0 に字句を要求し、取得した字句を字句  $t$  とする。字句  $t$  が入力文書の終端であった場合は、生成された抽象構文木を出力して処理を終了する（ステップ 2 0 3）。

#### 【 0 0 4 1 】

字句  $t$  が入力文書の終端でない場合、構文解析部 2 0 は、字句  $t$  が文脈ポインタに対して文法的に適合するかどうかを文法情報オブジェクト 2 4 に問い合わせる。そして、適合するならば文脈ポインタに字句  $t$  を追加する（ステップ 2 0 4 で Yes）。この追加は、次の処理によって行う。まず、非終端記号のノード  $n$  をノード作成部 2 3 により作成し、文脈ポインタに追加する（ステップ 2 0 5）。次に、文脈ポインタの指示先を新たに追加した非終端記号のノード  $n$  にずらす（ステップ 2 0 6）。そして、文脈ポインタが指し示す非終端記号のノード  $n$  が全ての子供のノードを持った場合、文脈ポインタを親ノードへずらす（ステップ 2 0 7、2 0 8）。文脈ポインタが指し示す非終端記号のノード  $n$  が全ての子供のノードを持っていない場合、またはステップ 2 0 8 で文脈ポインタを親ノード

へずらした後、ステップ 2 0 2 へ戻り、次の字句を取得して同様の処理を行う。

【 0 0 4 2 】

ステップ 2 0 4 において、字句  $t$  が文脈ポインタに対して文法的に適合しない場合、構文解析部 2 0 は、エラーの内容を出力し（ステップ 2 0 9）、構文回復部 2 2 に当該字句  $t$  における構文規則の誤りの修正を依頼する。字句  $t$  における構文規則の誤りを修正する処理は、ステップ 2 0 1 の構文解析部 2 0 の初期化処理において設定された構文回復部 2 2 を順次適用して誤りの修正を試行することにより行われる。構文回復部 2 2 を適用する順番は任意に設定可能であり、例えば適当な基準にて決定された優先度に基づいて定めることができる。

まず、未だ実行していない構文回復部 2 2 があるかどうかを調べ（ステップ 2 1 0）、未実行の構文回復部 2 2 がある場合は、そのうちの一つを実行する（ステップ 2 1 1）。そして、誤りの修正に成功した場合は、ステップ 2 0 2 へ戻り、次の字句に対する処理を行う（ステップ 2 1 2）。

【 0 0 4 3 】

誤りの修正に失敗した場合、ステップ 2 1 0 に戻り、未実行の構文回復部 2 2 を選んで実行する。ステップ 2 1 0 において、未実行の構文回復部 2 2 が存在しない場合、当該入力文書に対応するとして用意された構文回復部 2 2 によっては当該字句  $t$  における構文規則の誤りを修正することができなかったので、予め定められたエラー処理を実行した後（ステップ 2 1 3）、ステップ 2 0 2 へ戻る。ここで、エラー処理とは、当該字句  $t$  を読み飛ばして次の字句に対する処理に移行することや、構文解析部 2 0 自身が文脈ポインタに適当な字句  $t$  を付与するといった処理である。

【 0 0 4 4 】

上記の動作において、各構文回復部 2 2 は、次のようにして誤りの修正を実行する。

すなわち、まず、構文解析部 2 0 から、ステップ 2 0 4 の解析により得られた誤りの種類及び処理対象である字句  $t$  を付加情報として受け取る。この付加情報に基づいて、修正可能かどうかを判断する。修正不可能と判断したならば、次の構文回復部 2 2 に処理を渡し、次の構文回復部 2 2 が存在しないときは修正不可

能であることを構文解析部 2 0 に通知する。

一方、修正可能である場合は、それまでに生成されている抽象構文木、構文解析部 2 0 のバッファ 2 1 及び文脈ポインタを設定し直すことにより、当該字句 t に関する箇所を、構文規則の誤りがない状態に回復させる。そして、修正が成功したことを構文解析部 2 0 に通知する。

#### 【 0 0 4 5 】

図 3 は、本実施の形態の字句解析部 1 0 による字句解析処理の流れを説明するフローチャートである。

図 3 に示すように、字句解析処理が開始されると、まず、字句解析部 1 0 が初期化される（ステップ 3 0 1）。字句解析部 1 0 の初期化は、バッファ 1 1 を初期化する作業により行われる。また、入力文書の種類により使用する字句回復部 1 2 を選択する場合は、この初期化の時点で、対応表 3 0 などを参照して使用する字句回復部 1 2 を設定する。

#### 【 0 0 4 6 】

次に、字句解析部 1 0 は、構文解析部 2 0 から字句を要求されるのを待って、バッファ 1 1 から文字列 str を取り出して字句化し、構文解析部 2 0 に渡す（ステップ 3 0 2、3 0 3、3 0 4）。バッファ 1 1 が空の場合は（ステップ 3 0 1 で初期化した直後は常に空の状態となっている）、字句解析部 1 0 が入力文書の入力ストリームから文字列 str を読み込み、解析して構文解析部 2 0 に渡す。字句を構文解析部 2 0 に渡した後、字句解析部 1 0 はステップ 3 0 2 に戻り、構文解析部 2 0 から字句を要求されるのを待って、次の文字列 str に対する処理を行う。

#### 【 0 0 4 7 】

ステップ 3 0 3 において、文字列 str を字句化できなかった場合、字句回復部 1 2 を順次適用して修正処理を行う（ステップ 3 0 5、3 0 6、3 0 7）。字句回復部 1 2 は、文字列 str を正しい文字列に置き換え、バッファ 1 1 に挿入し、修正が成功したことを字句解析部 1 0 に通知する。そして、字句解析部 1 0 はステップ 3 0 2 に戻り、構文解析部 2 0 から字句を要求されるのを待って、次の文字列 str に対する処理を行う。

## 【 0 0 4 8 】

全ての字句回復部 1 2 によっても文字列strの修正が成功しなかった場合、最後の字句回復部 1 2 が、修正が成功しなかったことを字句解析部 1 0 に通知し、字句解析部 1 0 は、予め定められたエラー処理を実行した後（ステップ 3 0 8）、ステップ 3 0 2 へ戻る。ここで、エラー処理とは、当該文字列strを無視して次の文字列に対する処理に移行することなどである。

## 【 0 0 4 9 】

次に、本実施の形態による具体的な誤りの修正例について説明する。

ここでは、本実施の形態による具体的なシステムとして、ウェブページ上のコンテンツをオブジェクトとして扱うためのDOM (Document Object Model) を生成するSGMLパーザー (Parzer : 構文解析装置) に基づいたHTMLパーザーを取り上げる。SGMLのようなマークアップ言語では、構文解析システム 1 0 0 によって検出されるエラー (構文規則上の誤り) には、以下のようなものがある。ここでは、これを構文解析部 2 0 が構文回復部 2 2 に渡す誤りの種類とする。

## 【 0 0 5 0 】

## 1. 未定義属性

例えば、

```
<BODY background="bg.gif">
```

という開始タグがあった場合、background属性はHTML 4.0 Strictでは定義されていないため、エラーとなる (HTML 3.2 や HTML 4.0 Transitional には定義されている)。

## 【 0 0 5 1 】

## 2. 非対応終了タグ

(1) 本来あるべき位置よりも前方にある場合

例えば、

```
<p>str0 <b> str1 <i> str2 </b> str3 </i> str4 </p>
```

という断片で、入れ子構造をなすためには、</b>は</i>よりも後方にあるべきであるため、エラーとなる。

## (2) 対応すべき開始タグがない場合

例えば、

```
<table>
<td>..

```

という断片において、</tr>に対する開始タグ<tr>がないため、エラーとなる。

また、

```
<p>str0 <b> str1 <i> str2 </b> str3 </i> str4 </p>
```

という断片において、</i>に対応する開始タグがあった場合でも、</b>を取り扱う構文回復部 2 2 が</b>の直前に</i>があると解釈した場合、str3の直後の</i>に対応する開始タグが無いこととなるため、エラーとなる。

【 0 0 5 2 】

### 3. 不正エレメント（文脈ポインタの子供として文法的に正しくないエレメント）

例えば、

```
<table>
<form>
<tr><td>..

```

という断片において、TABLEの直下にはFORMは入らないという規則があるため、エラーとなる。

また、

```
<a name="top">
<p>
. . . . .
<a href="page1.html"> AnotherPage </a>
. . . . .
```



</a>

という断片において、アンカーエレメント（タグ<a>）の下にあるアンカーエレメントは例外として認められないという規則があるため、エラーとなる。

【 0 0 5 3 】

#### 4. 未定義エレメント

例えば、FONTエレメントはHTML 4.0 Strictでは定義されていないため、エラーとなる（HTML 3.2やHTML 4.0 Transitional には定義されている）。

【 0 0 5 4 】

次に、構文規則の誤りを含む入力文書の例を挙げて、誤りを修正する動作を説明する。

ここでは、誤りとして、上述した2. 非対応終了タグの（1）本来あるべき位置よりも前方にある場合の誤りと、（2）対応すべき開始タグがない場合の誤りとをそれぞれ含む、下記の文書Aと文書Bを例として取り上げる。なお、文書A、Bにおいて、誤り2.（1）発生、誤り2.（2）発生とコメントのある箇所が、それぞれ上述した誤りを含む箇所である。

【 0 0 5 5 】

文書A：

```
<table>
  <tr>
    <td>
      <form action="/invention1.html">
        <input name="input1" size=5>
      </td>      <!-- 誤り2.（1）発生 -->
    <td>
      <input type=submit>
    </form>      <!-- 誤り2.（2）発生 -->
  </td>
</tr>
```

</table>

【 0 0 5 6 】

文書 B :

<table>

<tr><td><form action="./invention1.html"><input name="input1" size  
=5><input type=submit></td></tr>

<tr>

<form>

<td>

<select>

<option>

. . . . .

<option>

</select>

<td>

</form>

</tr>

</table>

<table>

<tr>

<td>

<select>

<option>

. . . . .

<option>

</select>

</td>

</tr>

</table>

</form> <!-- 誤り 2. (2) 発生 -->

【 0 0 5 7 】

図 4 は、上記文書 A に対して構文回復部 2 2 による誤りの修正を行わずに構文解析部 2 0 による構文解析を行った場合における出力結果の抽象構文木を示す図である。

図 4 を参照すると、2 番目の INPUT エLEMENT が FORM エLEMENT の下に位置していないという誤りが発生している。これは、文書 A における 1 番目の </td> タグが出現した場合にそれに対応する TD エLEMENT を閉じた結果、FORM エLEMENT も同時に閉じてしまうという構文回復部 2 2 による修正なしの場合のデフォルト処理が原因となっている。

【 0 0 5 8 】

図 5 は、上記文書 B に対して構文回復部 2 2 による誤りの修正を行わずに構文解析部 2 0 による構文解析を行った場合における出力結果の抽象構文木を示す図である。

図 5 を参照すると、これ 2 番目の SELECT エLEMENT が FORM の下に位置していないという誤りが発生している。これは、文書 B における最後の </form> タグに対応する開始タグ <form> が存在しないため、</form> タグを無視して処理を終了したことが原因となっている。

【 0 0 5 9 】

これらの FORM タグに関する誤りに対し、構文回復部 2 2 として、FormExpander と FormInserter とを用いて誤りを修正する場合を考える。FormExpander 及び FormInserter の起動条件と誤りを修正する動作（回復動作）は以下に示す通りである。

【 0 0 6 0 】

[FormExpander]

起動条件：誤りを起こした字句  $t$  が FORM エLEMENT の終了タグ </form> であり、現在の文脈ポインタが指し示すノードからルート方向に辿っても FORM エLEMENT が見つからず、かつ当該字句  $t$  が文脈ポインタの前方にある場合。

回復動作：

1. 文脈ポインタから前方へ走査し、FORMエレメントfe0を探す。
2. FORMエレメントfe0を元々あった位置p0から除く。ここで言う「除く」とは、当該ノードが保持している子孫のノードは削除せずにFORMエレメントfe0のみを取り除くことであり、当該子孫のノードはFORMエレメントfe0の親ノードに追加する。
3. 文脈ポインタと位置p0とを最小の範囲で覆うことができるような位置にFORMエレメントfe0を挿入する。
4. 文脈ポインタがFORMエレメントfe0の親ノードを指すようにする。

【 0 0 6 1 】

[FormInserter]

起動条件：誤りを起こした字句tがFORMエレメントの終了タグ</form>であり、現在の文脈ポインタが指し示すノードからルート方向に辿ってもFORMエレメントが見つからず、かつ当該字句tが文脈ポインタの前方にある場合。

回復動作：

1. 文脈ポインタの前方へFORMエレメントfe1を探す。
2. FORMエレメントfe2を作る。
3. 文脈ポインタの指し示すノードと共通の親ノードを持つ他のノードのなかで、FORMエレメントfe1よりも後方にあるものを、文脈ポインタの指し示すノードの親ノードcp0から切り取り、FORMエレメントfe2に子ノードとして追加する。ただし、ここで言う「切り取る」とは、上述した「除く」と違い、当該ノードが保持している子孫のノードも含めた部分木を取り除くことである。
4. FORMエレメントfe2をノードcp0に追加する。
5. 文脈ポインタがノードcp0を指すようにする。

【 0 0 6 2 】

上述したように、FormExpanderとFormInserterとは、起動条件が同一である。すなわち、構文解析において、入力文書から上記の起動条件に合致する構文規則の誤りが検出された場合、FormExpanderまたはFormInserterのいずれか一方を選択的に用いることにより、当該誤りを修正することが可能となる。

【 0 0 6 3 】

図 6 は、文書 A に対して FormExpander を用いて誤りを修正した場合における構文解析部 2 0 による構文解析の出力結果の抽象構文木を示す図である。

図 6 を参照すると、TR エLEMENT の下に FORM エLEMENT が来るという問題はあるが、二つの INPUT エLEMENT がそれぞれ FORM エLEMENT の下に位置していることがわかる。

#### 【 0 0 6 4 】

図 7 は、文書 B に対して FormExpander を用いて誤りを修正した場合における構文解析部 2 0 による構文解析の出力結果の抽象構文木を示す図である。

図 7 を参照すると、2 番目の SELECT エLEMENT は FORM の下に位置させることができた。しかし、FORM の下に FORM が位置する例外に該当してしまっている。

#### 【 0 0 6 5 】

図 8 は、文書 A に対して FormInserter を用いて誤りを修正した場合における構文解析部 2 0 による構文解析の出力結果の抽象構文木を示す図である。

図 8 を参照すると、本来は同じ FORM エLEMENT の下に位置しなければならない二つの INPUT エLEMENT が、別々の FORM エLEMENT の下に位置している。

#### 【 0 0 6 6 】

図 9 は、文書 B に対して FormInserter を用いて誤りを修正した場合における構文解析部 2 0 による構文解析の出力結果の抽象構文木を示す図である。

図 9 を参照すると、2 番目の SELECT エLEMENT が FORM の下に位置しており、かつ FORM エLEMENT が FORM エLEMENT を含んでいないことがわかる。

#### 【 0 0 6 7 】

以上の FormExpander 及び FormInserter の適用結果を参酌すると、文書 A に示す誤りを修正するには FormExpander が適しており、文書 B に示す誤りを修正するには FormInserter が誤り回復器として適していることがわかる。したがって、文書 A を構文解析する場合と、文書 B を構文解析する場合とで、使用する構文回復部 2 2 を変更することが好ましい。

#### 【 0 0 6 8 】

図 1 0 は、上記の動作例において、適切な構文回復部 2 2 を適用するための対応表 3 0 の構成例を示す図表である。

図 1 0 を参照すると、文書の識別子として「文書 A」、「文書 B」が登録され（実際には、文書 A、B を識別する所定の ID が登録される）、「文書 A」に対応する構文回復部 2 2 として FormExpander が登録され、「文書 B」に対応する構文回復部 2 2 として FormInsertter が登録されている。

これにより、文書 A を構文解析する場合には、構文解析部 2 0 を初期化した際に、文書 A の ID をキーとして構文回復部 2 2 として FormExpander を使用するよう設定される。同様に、文書 B を構文解析する場合には、構文解析部 2 0 を初期化した際に、文書 B の ID をキーとして構文回復部 2 2 として FormInsertter を使用するよう設定される。そして、適切な構文回復部 2 2 を用いて入力文書（文書 A または文書 B）の誤りを修正することができ、最適な出力結果（抽象構文木）を得ることができる。

#### 【0069】

以上、FORM タグに関する誤りを修正する構文回復部 2 2 について説明したが、次に、他のいくつかの構文回復部 2 2 について例を挙げて説明する。なお、以下に示す構文回復部 2 2 は、HTML 文書に対する構文解析に用いられる構文回復部 2 2 の例であるが、HTML 文書以外の SGML 文書や XML 文書にそのまま適用したり、一般化して適用したりできるものもある。

#### 【0070】

##### [DefaultErrorHandler]

起動条件：誤りの種類が不正エレメント（上述したエラーの種類の 3）であり、誤りを起こしたエレメントの正しい出現場所が固定的である場合。ここで、固定的とは、DTD（Document Type Definition：文書型定義）によって、エレメントの位置が繰り返しのない親エレメントの下に固定されていることを意味する。例えば、HTML エレメントの下 HEAD エレメントなどが該当する。

回復動作：

1. 不正エレメント e1 と同名のエレメント e2 を探す。
2. エレメント e2 が他の構文回復部 2 2 の処理または構文解析部 2 0 の処理によって補完されたエレメントかどうかを調べる。
3. エレメント e2 が補完されたエレメントである場合、エレメント e2 をエレメン

トe1で置き換える。

この構文回復部 2 2 は、HTMLとは独立であるので、SGML文書やXML文書に適用することができる。

【 0 0 7 1 】

[InterleavedEndtagExchanger]

起動条件：誤りの種類が非対応終了タグであり、当該タグが本来あるべき位置よりも前方にある場合（上述したエラーの種類の 2.（1））。

回復動作：

1. エレメントを示す変数e1に文脈ポインタが指し示すノードの親ノードを代入する。
2. 変数e1が、誤りを起こした非対応終了タグet0と同じである場合、修正は不可能と判断し、処理を終了する（回復失敗）。
3. 変数e1が非対応終了タグet0と同じでなく、終了タグを省略不可能である場合、構文解析部 2 0 のバッファ 2 1 から後方に向かって、変数e1と同名の終了タグet1を探索する。
4. ステップ 3 の探索により終了タグet1を発見できた場合、構文解析部 2 0 のバッファ 2 1 におけるエレメントet1が入っていた位置に、非対応終了タグet0を入れ、バッファ 2 1 の先頭に終了タグet1を入れる。
5. ステップ 3 の探索で終了タグet1を発見できなかった場合、変数e1を当該変数e1の親ノードに変更してステップ 2 へ戻る。

この構文回復部 2 2 は、HTMLとは独立であるので、SGML文書やXML文書に適用することができる。

【 0 0 7 2 】

[RangeExpander]

起動条件：誤りの種類が非対応終了タグであり、当該タグが本来あるべき位置よりも前方にある場合（上述したエラーの種類の 2.（1））。

回復動作：

1. エレメントを示す変数e1に文脈ポインタを代入し、エレメントの配列型の変数配列を作成する。

2. 変数e1が誤りを起こした非対応終了タグet0と同じ名前かどうかを判断し、同じでなければ同じ名前になるまで、以下のステップ3、4を繰り返す。
3. 変数e1が終了タグを省略不可能である場合、ノード作成部23を用いて変数e1と同名のエレメントを作成し、ステップ1で作成した変数配列の末尾に挿入する。
4. 変数e1をエレメントe1にずらしてステップ2へ戻る。
5. ステップ2において、変数e1が誤りを起こした非対応終了タグet0と同じ名前であるならば、変数配列が空かどうかを判断し、空でなければ空になるまでステップ6を繰り返す。
6. 変数配列の先頭からエレメントe2を取り出し、文脈ポインタにエレメントe2を追加する。そして、文脈ポインタの指示先をエレメントe2に変更する。

この構文回復部22は、HTMLとは独立であるので、SGML文書やXML文書に適用することができる。

#### 【0073】

##### [AnchorUnderAnchorHandler]

起動条件：誤りの種類が不正エレメント（上述したエラーの種類の3）であり、誤りを起こしたエレメントa1がアンカーであり、文脈ポインタが指し示すノードの上方にもアンカーであるエレメントa2がある場合。

回復動作：エレメントa1を、エレメントa2と共通の親ノードを持つノードとして追加し、文脈ポインタの指示先をエレメントa1に変更する。

この構文回復部22は、上方エレメントueの下で下方エレメントdeが例外となる誤りの場合に、上方エレメントdeを下方エレメントueと共通の親を持つノードにするという回復方法として一般化できる。これにより、SGML文書やXML文書にも適用することができる。

#### 【0074】

##### [FramesetErrorHandler]

起動条件：誤りの種類が未定義エレメント（上述したエラーの種類の4）であり、誤りを起こしたエレメントe1がFRAMESETである場合。

回復動作：



1. エレメントe1を構文解析部20のバッファ21に戻す。
2. 構文解析部20が使用している文法情報オブジェクト24を、FRAMESETを定義しているHTML4.0 Framesetなどに変更する。

この構文回復部22は、未定義エレメントe2に対し、エレメントe2を定義している文書型定義DTD2に文法情報オブジェクト24を設定する方法として一般化できる。これにより、SGML文書やXML文書にも適用することができる。文書型定義DTD2は、使用者により明示的に指定しても良いし、DTDのレポジトリなどから探索しても良い。

【0075】

[HTMLErrorHandler]

起動条件：誤りの種類が不正エレメント（上述したエラーの種類の3）であり、誤りを起こしたエレメントe1の名前が

LINK, STYLE, META, BASE, ISINDEX

のいずれかである場合。

回復動作：出力結果ASTのルートノードにあるHTMLエレメントの子ノードからHEADエレメントを探し、当該HEADエレメントのノードにエレメントe1を追加する。

この構文回復部22は、HTMLに固有のものである。

【0076】

[IgnoreFont]

起動条件：誤りの種類が不正エレメント（上述したエラーの種類の3）であり、誤りを起こしたエレメントe1の名前がFONTである場合。

回復動作：何ら処理を行わず、エレメントe1を無視する。

この構文回復部22は、一般化してSGML文書やXML文書にも適用することができる。

【0077】

[TRErrorHandler]

起動条件：誤りの種類が不正エレメント（上述したエラーの種類の3）であり、誤りを起こしたエレメントe1の名前がTDである場合。

回復動作：文脈ポインタの指示先がTBODYエレメントまたはTABLEエレメントであ

る場合に、新たにTRエレメントtrを作成し文脈ポインタに追加する。そして、エレメントtrに対してエレメントelを追加し、文脈ポインタの指示先をエレメントelに設定する。

#### 【0078】

次に、字句回復部12の例を説明する。字句解析においては、文字列が予め定められたタグの字句などに該当するかどうかを調べるだけなので、字句回復部12の種類は多くはない。以下に、代表的な例を挙げる。

[AttributeValueErrorHandler]

起動条件：開始タグ中の属性値を字句解析中に、

"... > ... < ..."

という文字列が出現した場合。

回復動作："... "> ... < ... "という文字列を字句解析部10のバッファ11に挿入する。

#### 【0079】

なお、上記の説明では、処理対象となる入力文書がマークアップ言語で記述された文書、特にHTML文書である場合について説明したが、他のプログラミング言語で記述されたプログラムや自然言語で記述された文書に対する構文解析にも、当該プログラミング言語や自然言語における字句の規則及び構文規則に対応する回復手段を用意することにより、そのまま利用することが可能である。

#### 【0080】

さらに、これらのプログラムや文書に限らず、楽譜のような一定の規則に基づいて記述されたデータ列に対する構造解析及び誤りの修正を行う構造回復システムとして利用することもできる。

この場合、当該構造回復システムは、データ列の構造を解析し、規則上の誤りを検出する解析手段（本実施の形態における構文解析部20に相当する）と、この解析手段とは独立に設けられ、解析手段からの依頼に応じて、データ列から検出された規則上の誤りを修正する回復手段（本実施の形態における構文回復部22の集合に相当する）とを備え、かつこの回復手段は、特定の種類の誤りを修正する単純な機能を有する修正手段（本実施の形態における個々の構文回復部22

に相当する)の集合を含んで構成される。そして、データ列における規則上の誤りの種類に応じてこの修正手段を切り換えて使用することにより、このデータ列における種々の誤りを修正することとなる。

【 0 0 8 1 】

【発明の効果】

以上説明したように、本発明によれば、プログラムや文書などのデータ列における字句や構文規則の誤りの種類に応じて複数の回復手段を用意し、選択的に適用して誤りを修正することができる。

【図面の簡単な説明】

【図 1】 本発明の実施の形態における構文解析システムの全体構成を示す図である。

【図 2】 本実施の形態の構文解析部による構文解析処理の流れを説明するフローチャートである。

【図 3】 本実施の形態の字句解析部による字句解析処理の流れを説明するフローチャートである。

【図 4】 構文回復部による誤りの修正を行わずに構文解析部による構文解析を行った場合における出力結果の抽象構文木の例を示す図である。

【図 5】 構文回復部による誤りの修正を行わずに構文解析部による構文解析を行った場合における出力結果の抽象構文木の他の例を示す図である。

【図 6】 構文回復部を用いて誤りを修正した場合における構文解析部による構文解析の出力結果の抽象構文木の例を示す図である。

【図 7】 構文回復部を用いて誤りを修正した場合における構文解析部による構文解析の出力結果の抽象構文木の他の例を示す図である。

【図 8】 他の構文回復部を用いて誤りを修正した場合における構文解析部による構文解析の出力結果の抽象構文木の例を示す図である。

【図 9】 他の構文回復部を用いて誤りを修正した場合における構文解析部による構文解析の出力結果の抽象構文木の他の例を示す図である。

【図 1 0】 本実施の形態において適切な構文回復部を適用するために用いる対応表の構成例を示す図表である。

【図 1 1】 本実施の形態の構文解析システムを搭載したコンピュータ装置の構成例を説明する図である。

【図 1 2】 従来の構文解析システムの構成例を示す図である。

【図 1 3】 図 1 2 に示す従来の構文解析システムにおける構文解析処理を説明するフローチャートである。

【図 1 4】 従来の構文解析システムによる HTML 文書の修正例を示す図である。

【図 1 5】 従来の構文解析システムによる HTML 文書の他の修正例を示す図である。

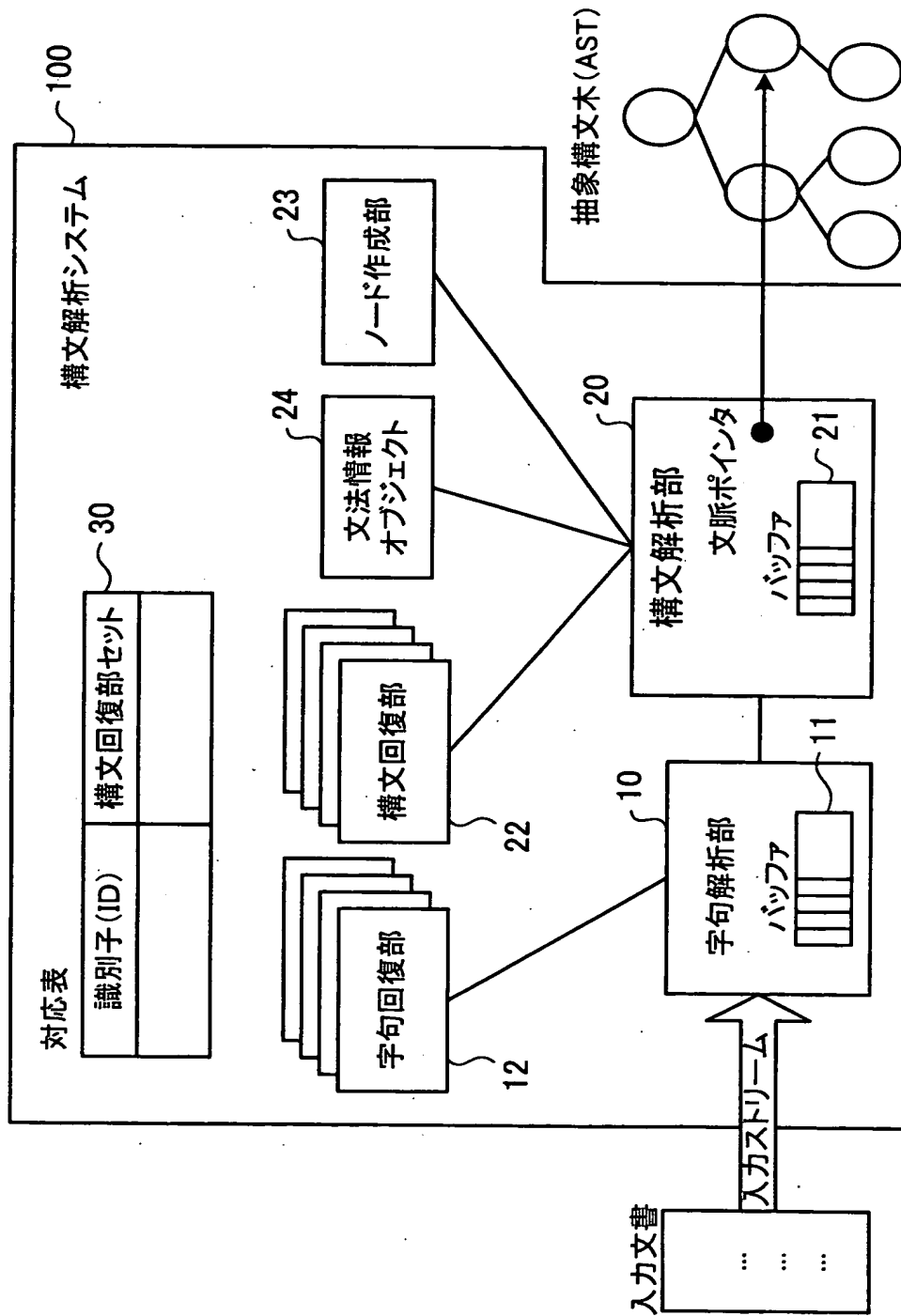
【符号の説明】

1 0 … 字句解析部、 1 1、 2 1 … バッファ、 1 2 … 字句回復部、 2 0 … 構文解析部、 2 2 … 構文回復部、 2 3 … ノード作成部、 2 4 … 文法情報オブジェクト、 3 0 … 対応表

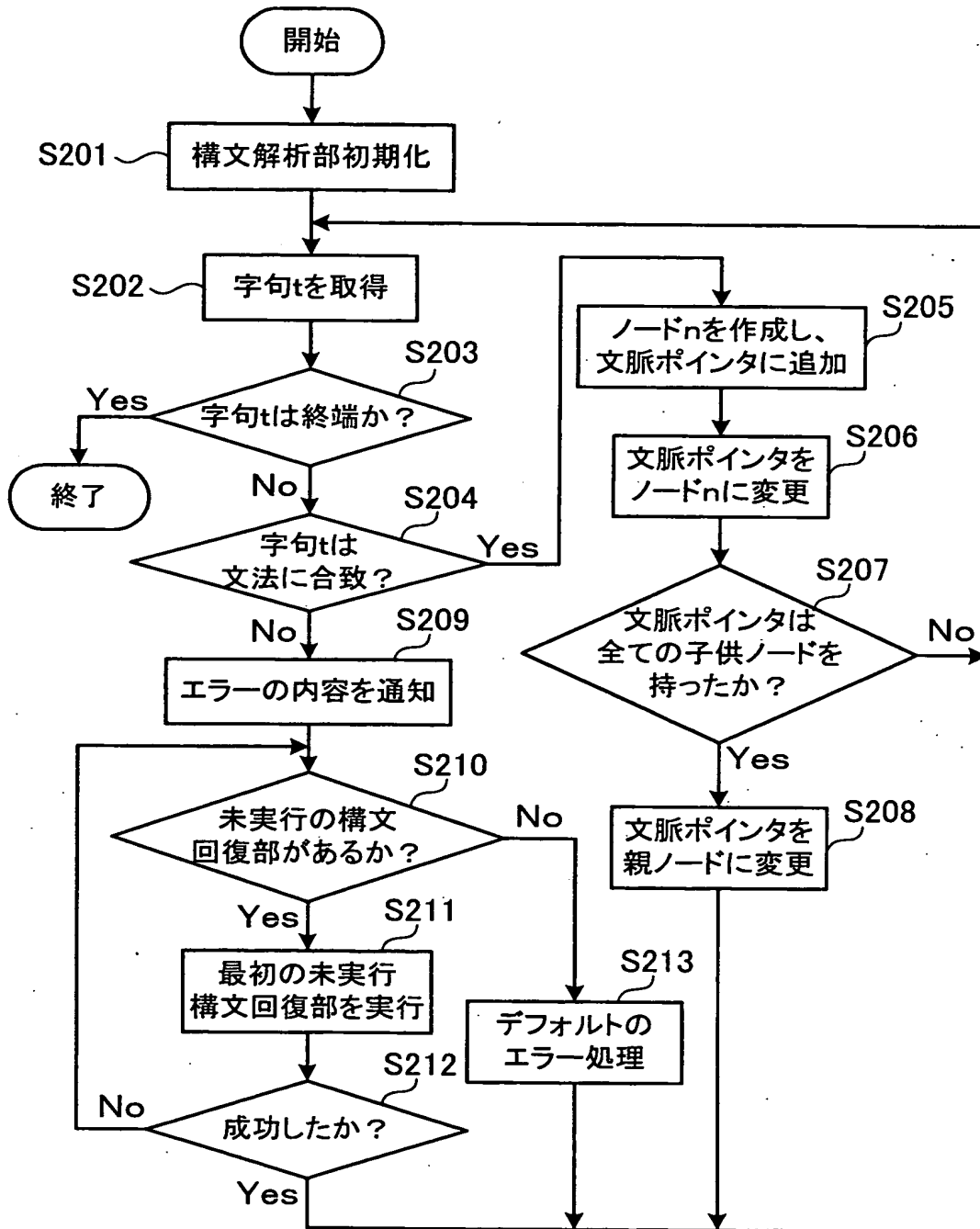
【書類名】

図面

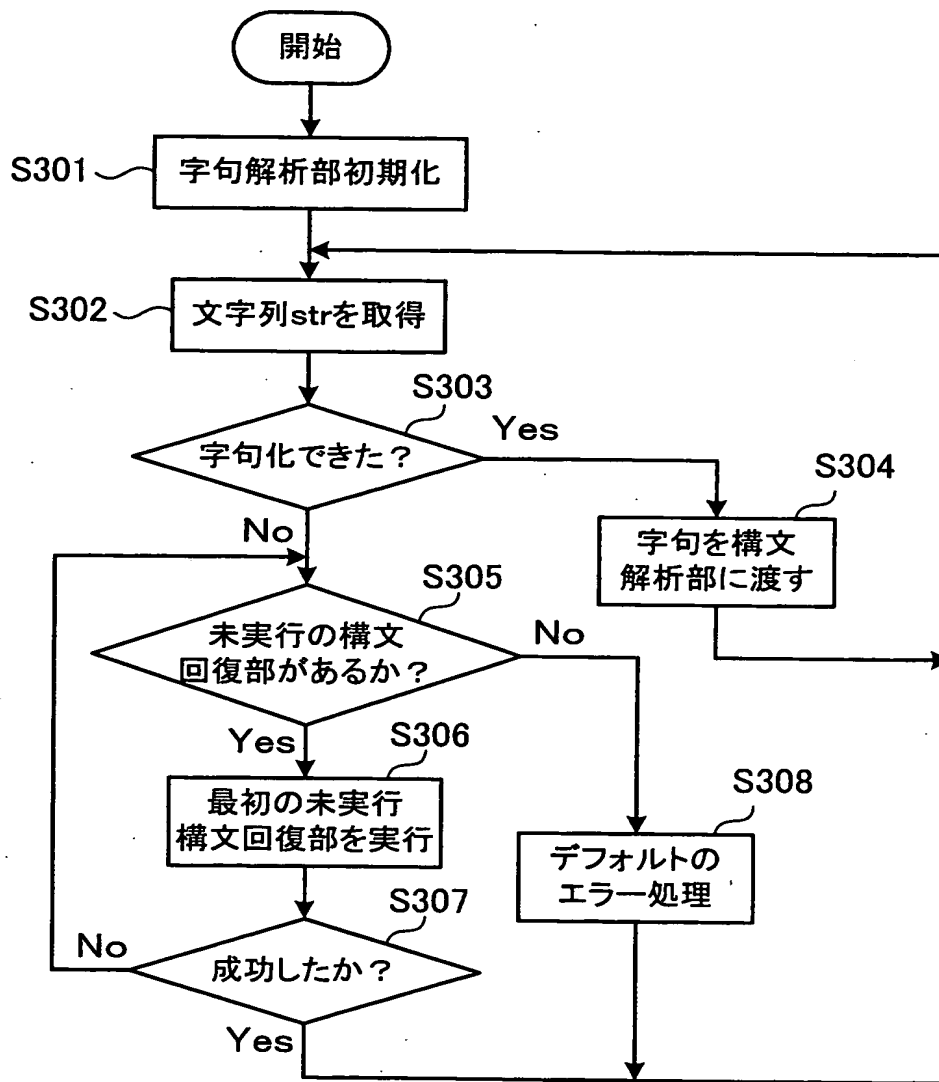
【図 1】



【図 2】

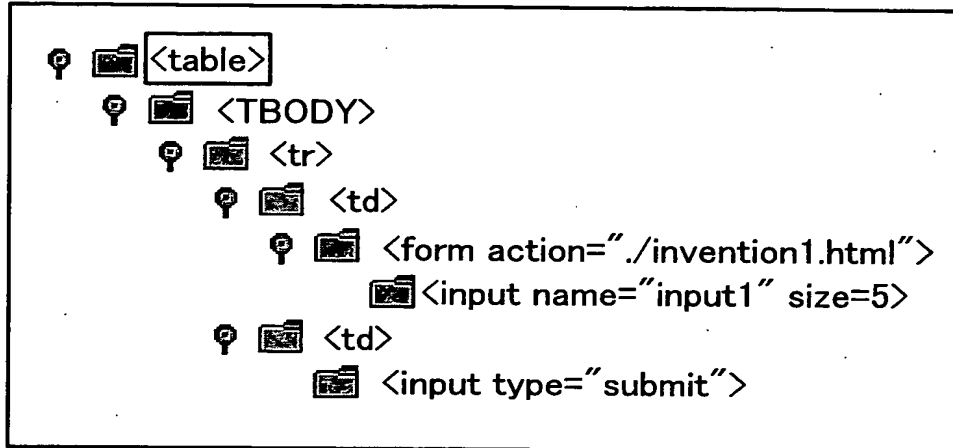


【図 3】



【図 4】

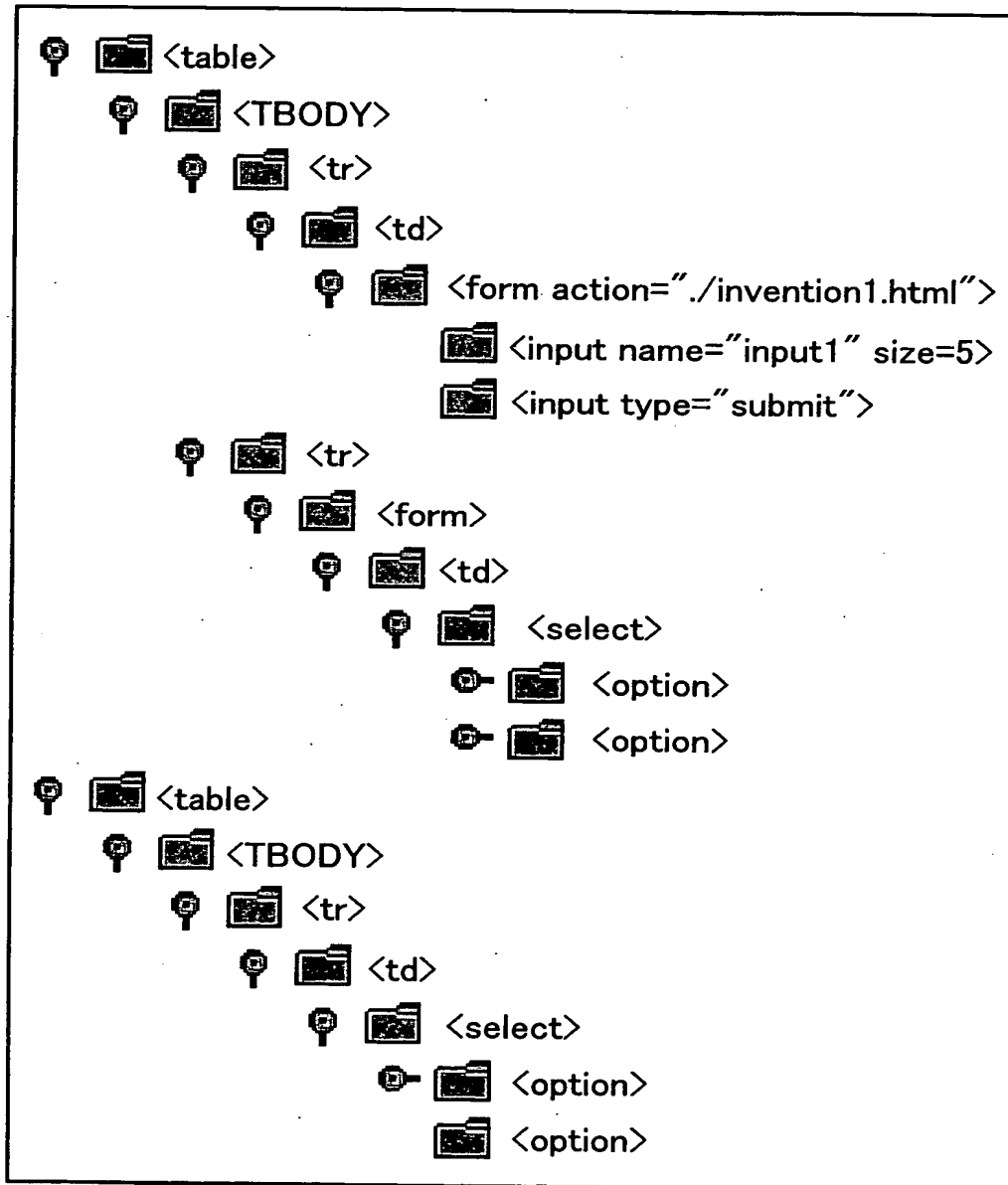
構文回復部による修正を行わない場合の文書Aの解析結果





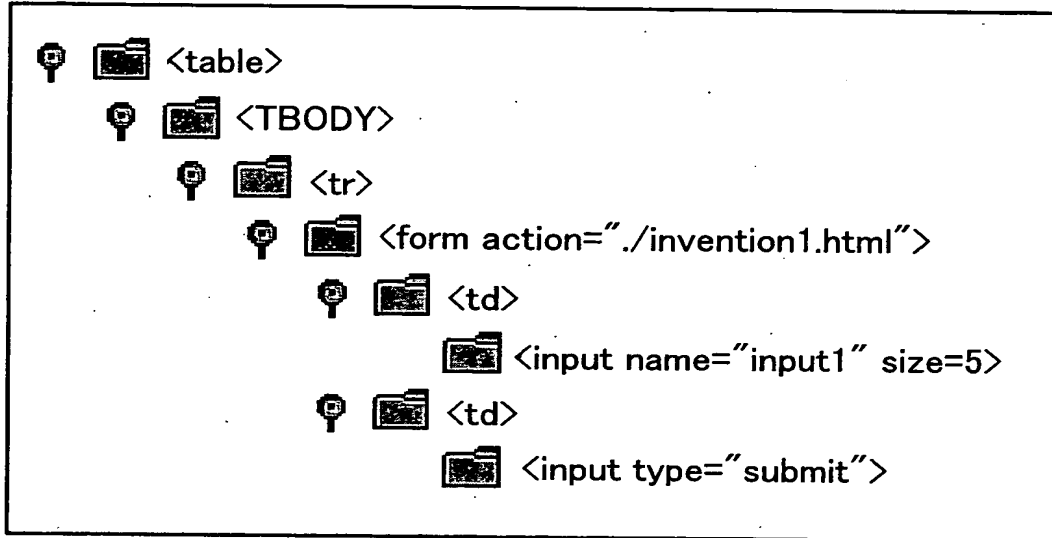
【図 5】

構文回復部による修正を行わない場合の文書Bの解析結果



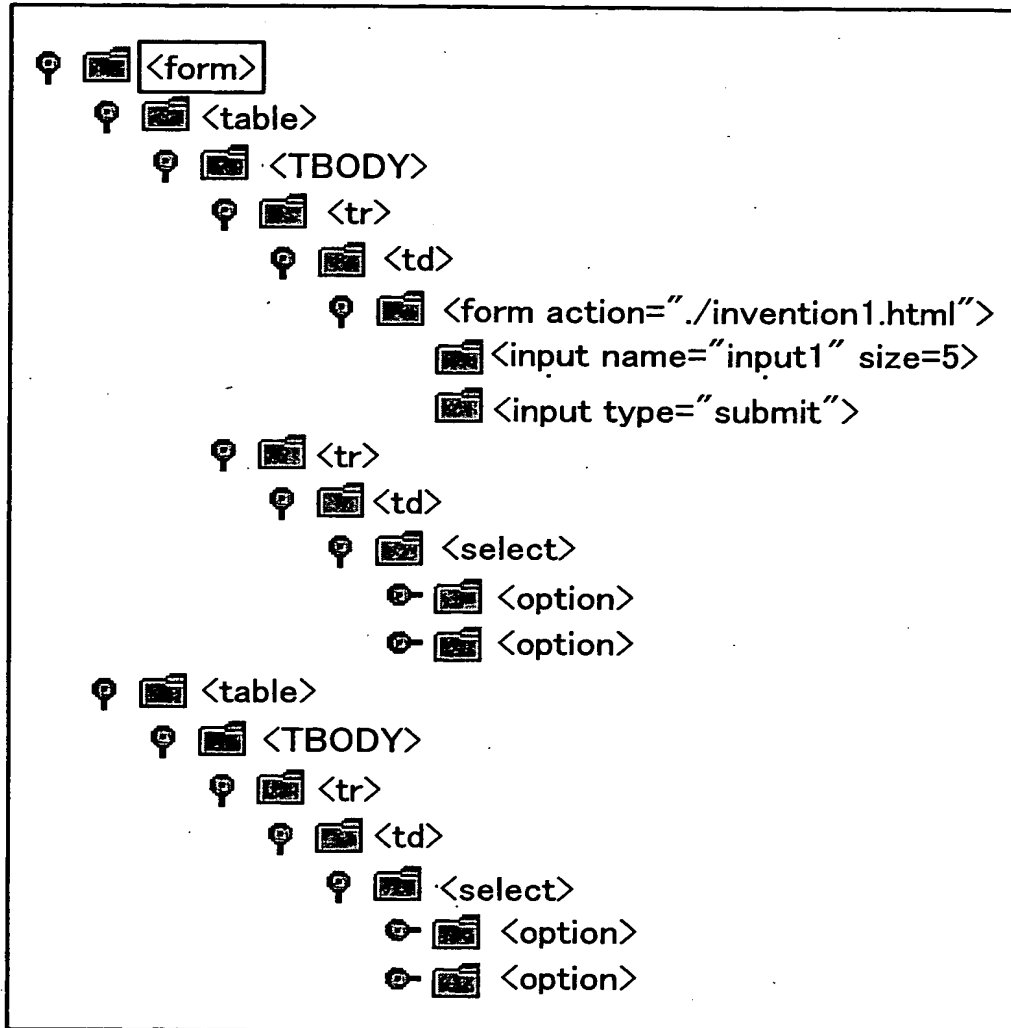
【図 6】

FormExpanderを適用した場合の文書Aの解析結果



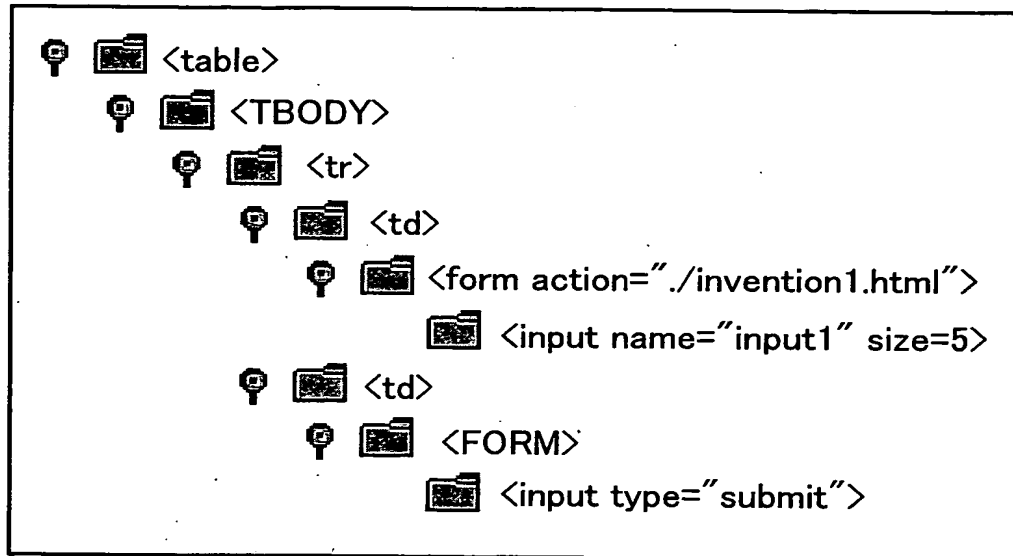
【図 7】

FormExpanderを適用した場合の文書Bの解析結果



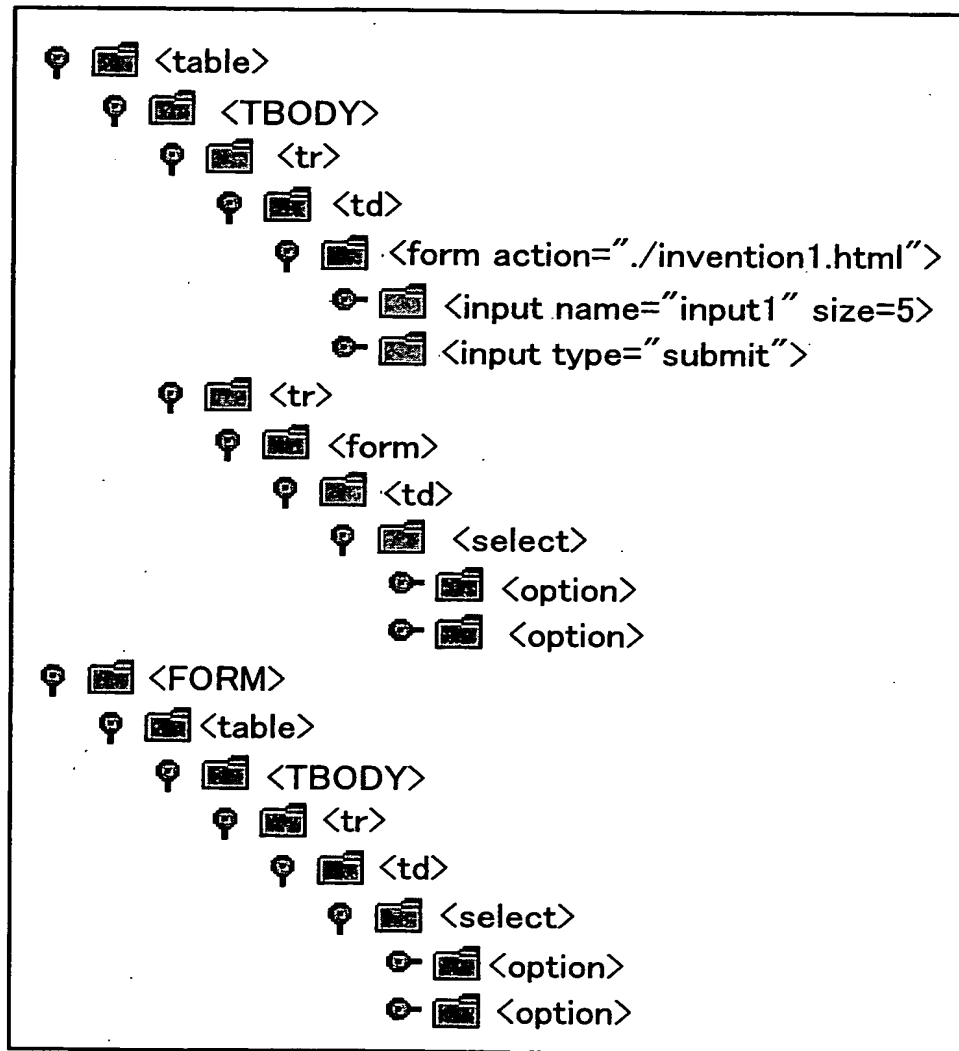
【図 8】

FormInserterを適用した場合の文書Aの解析結果



【図 9】

## FormInserterを適用した場合の文書Bの解析結果

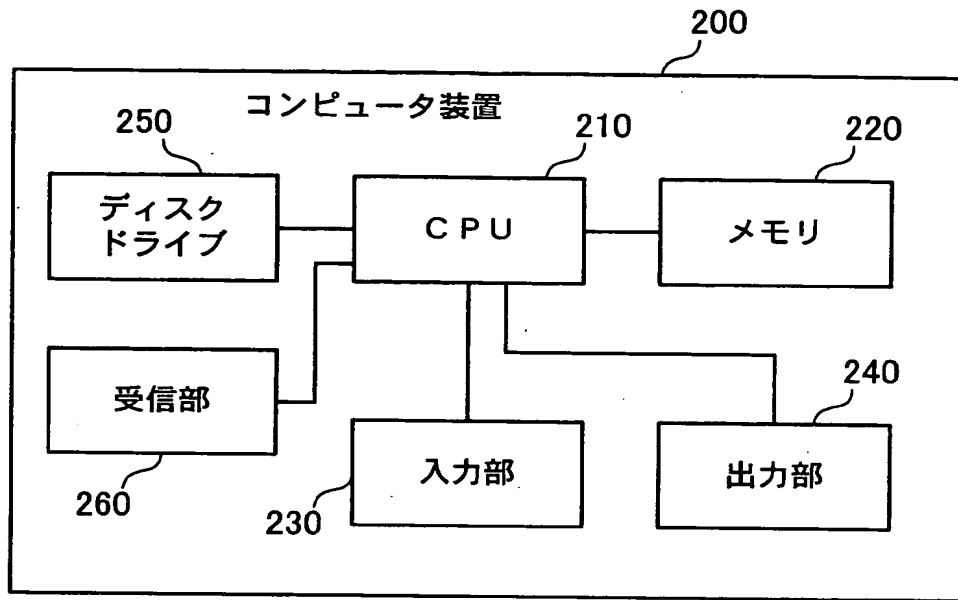


【図 1 0】

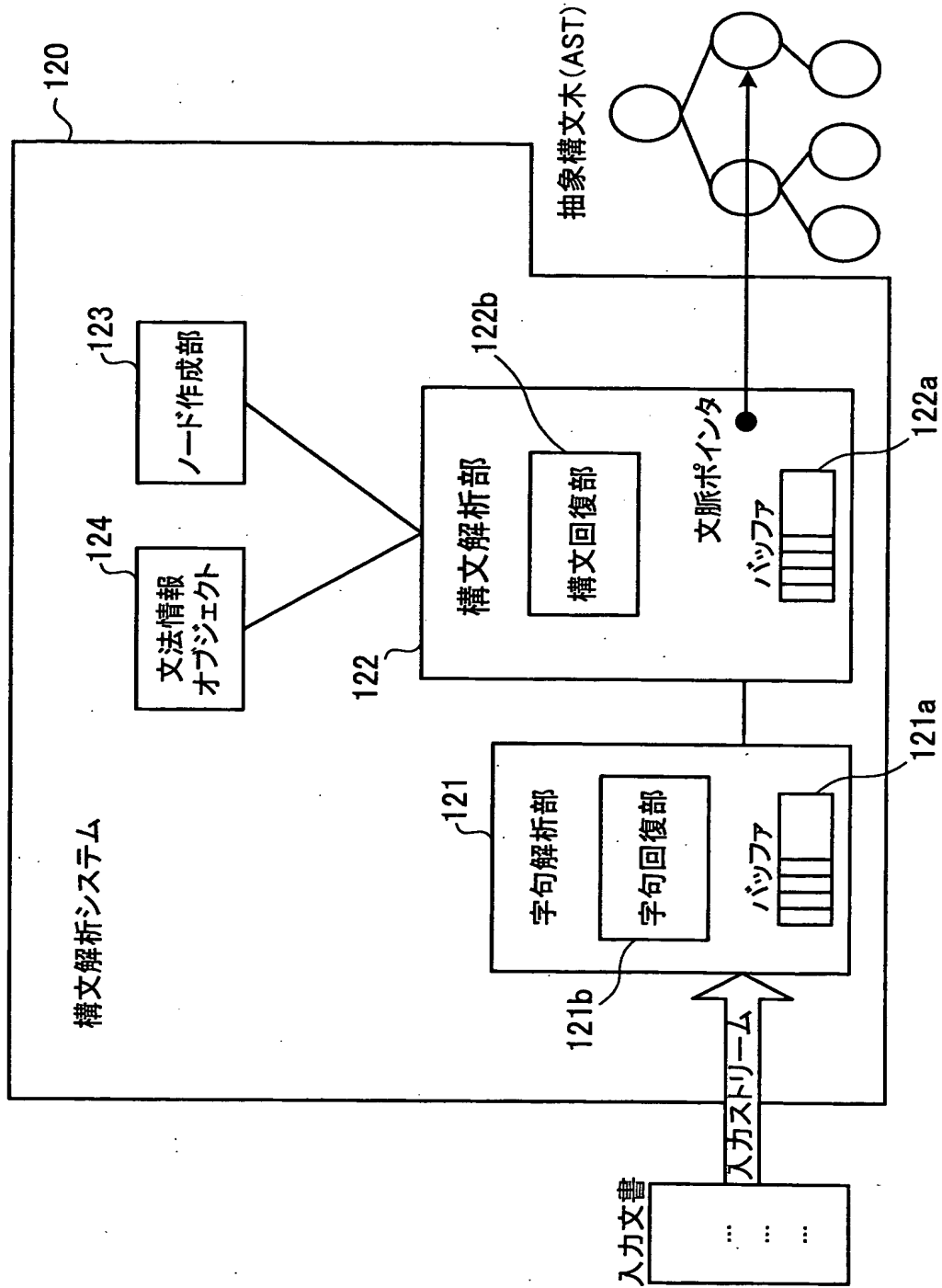
30

文書の識別子	構文回復部
文書 A	[FormExpander]
文書 B	[FormInserter]
その他	□

【図 1 1】

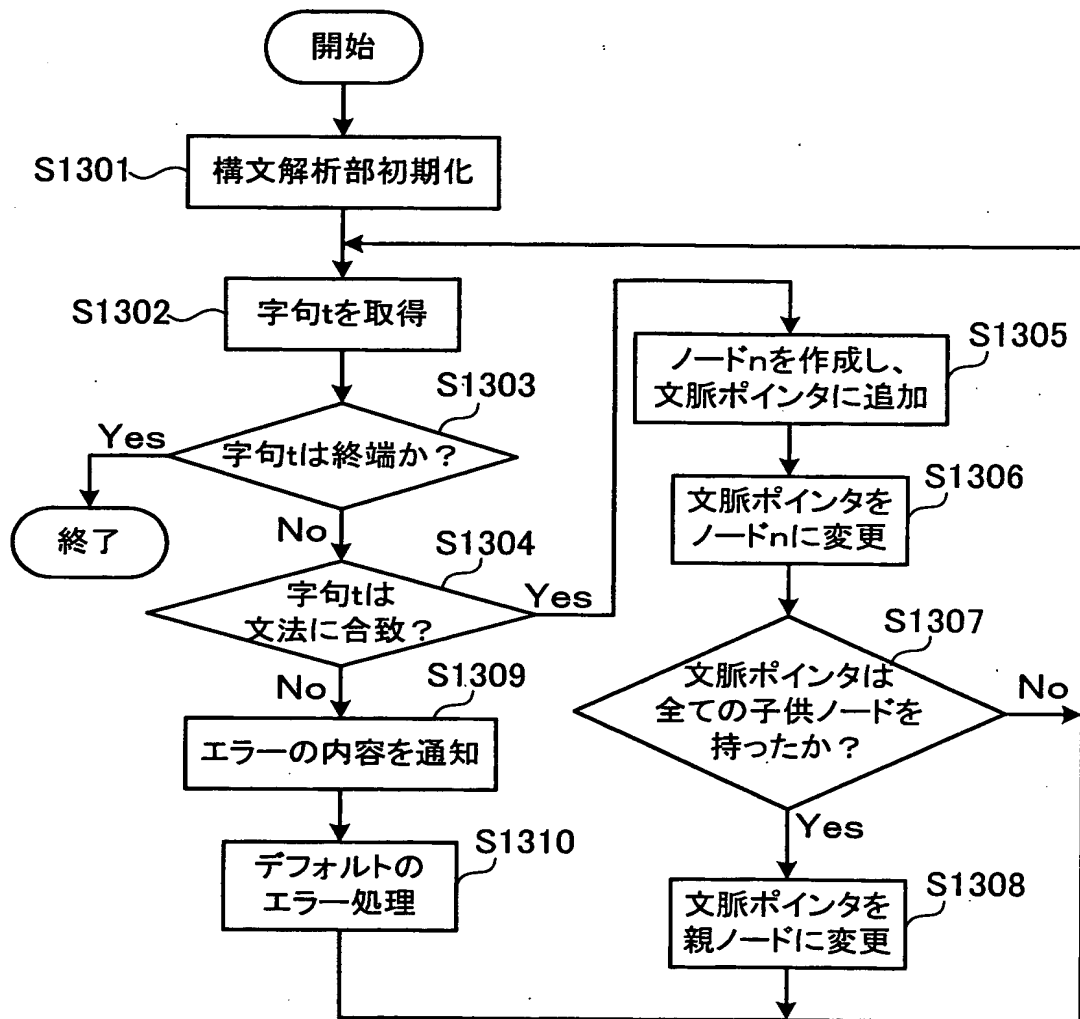


【図 12】





【図 1 3】



【図 1 4】

str0 str1 *str2* str3 str4

【図 1 5】

str0 str1 *str2 str3* str4

【書類名】 要約書

【要約】

【課題】 プログラムや文書などのデータ列における字句や構文規則の誤りの種類に応じて複数の回復手段を用意し、選択的に適用して誤りを修正できるシステムを提供する。

【解決手段】 所定の規則に基づいて記述されたデータ列の構造を解析し、規則上の誤りを検出する構文解析部 2 0 と、この構文解析部 2 0 とは独立に設けられ、この構文解析部 2 0 からの依頼に応じて、この構文解析部 2 0 により検出されたこのデータ列における規則上の誤りを修正する回復手段とを備え、この回復手段は、特定の種類の誤りを修正する単純な機能を有する構文回復部 2 2 の集合を含んで構成され、このデータ列における規則上の誤りの種類に応じてこの構文回復部 2 2 を切り換えて使用することにより、このデータ列における種々の誤りを修正する。

【選択図】 図 1

認定・付加情報

特許出願の番号	特願 2 0 0 0 - 3 2 3 8 2 3
受付番号	5 0 0 0 1 3 7 2 2 1 8
書類名	特許願
担当官	塩崎 博子 1 6 0 6
作成日	平成 1 2 年 1 2 月 8 日

<認定情報・付加情報>

【特許出願人】

【識別番号】	390009531
【住所又は居所】	アメリカ合衆国 1 0 5 0 4、ニューヨーク州 アーモンク (番地なし)
【氏名又は名称】	インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】	100086243
【住所又は居所】	神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	坂口 博

【代理人】

【識別番号】	100091568
【住所又は居所】	神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	市位 嘉宏

【代理人】

【識別番号】	100106699
【住所又は居所】	神奈川県大和市下鶴間 1 6 2 3 番 1 4 日本アイ・ビー・エム株式会社大和事業所内
【氏名又は名称】	渡部 弘道

【復代理人】

申請人	
【識別番号】	100104880
【住所又は居所】	東京都港区赤坂 5 - 4 - 1 1 山口建設第 2 ビル 6 F セリオ国際特許事務所
【氏名又は名称】	古部 次郎

【選任した復代理人】

【識別番号】	100100077
--------	-----------

次頁有

認定・付加情報（続き）

【住所又は居所】 東京都港区赤坂 5 - 4 - 1 1 山口建設第 2 ビル  
6 F セリオ国際特許事務所  
【氏名又は名称】 大場 充

出 願 人 履 歴 情 報

識別番号 [390009531]

1. 変更年月日 2000年 5月16日

[変更理由] 名称変更

住 所 アメリカ合衆国10504、ニューヨーク州 アーモンク (番地なし)

氏 名 インターナショナル・ビジネス・マシーンズ・コーポレーション